



## Stealth Attacks on the Smart Grid

Ke Sun, Iñaki Esnaola, Samir Perlaza, Harold Vincent Poor

### ► To cite this version:

Ke Sun, Iñaki Esnaola, Samir Perlaza, Harold Vincent Poor. Stealth Attacks on the Smart Grid. 2018. hal-01857366v1

**HAL Id: hal-01857366**

**<https://hal.inria.fr/hal-01857366v1>**

Preprint submitted on 15 Aug 2018 (v1), last revised 2 Feb 2021 (v2)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Stealth Attacks on the Smart Grid

Ke Sun, Iñaki Esnaola, Samir M. Perlaza, and H. Vincent Poor

## Abstract

Random attacks that jointly minimize the amount of information acquired by the operator about the state of the grid and the probability of attack detection are presented. The attacks minimize the information acquired by the operator by minimizing the mutual information between the observations and the state variables describing the grid. Simultaneously, the attacker aims to minimize the probability of attack detection by minimizing the Kullback-Leibler (KL) divergence between the distribution when the attack is present and the distribution under normal operation. The resulting cost function is the weighted sum of the mutual information and the KL divergence mentioned above. The trade-off between the probability of attack detection and the reduction of mutual information is governed by the weighting parameter on the KL divergence term in the cost function. The probability of attack detection is evaluated as a function of the weighting parameter. A sufficient condition on the weighting parameter is given for achieving an arbitrarily small probability of attack detection. The attack performance is numerically assessed on the IEEE 30-Bus and 118-Bus test systems.

## Index Terms

Stealth, data injection attacks, information-theoretic security, mutual information, probability of detection

## I. INTRODUCTION

The smart grid relies on the effective integration of the power grid and advanced communication and sensing infrastructure. Consistency between the physical layer of the power

K. Sun and I. Esnaola are with the Department of Automatic Control and Systems Engineering, University of Sheffield, Sheffield S1 3JD, UK. I. Esnaola is also with the Department of Electrical Engineering, Princeton University, Princeton NJ 08540, USA. (email: ke.sun@sheffield.ac.uk, esnaola@sheffield.ac.uk).

S. M. Perlaza is with the Institut National de Recherche en Informatique et Automatique (INRIA), Lyon, France, and also with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA (email: samir.perlaza@inria.fr).

H. V. Poor is with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: poor@princeton.edu).

grid and the energy management system (EMS) in the cyber layer facilitates an economic and reliable operation of the power system. The 2003 North American outage caused by an alarm system failure [1] and the 2015 Ukraine power failure caused by the BlackEnergy virus incident [2] emphasize the need for cybersecurity mechanisms for the power system. However, the cybersecurity threats to which the smart grid is exposed are not well understood yet, and therefore, practical security solutions need to come forth as a multidisciplinary effort combining technologies such as cryptography, machine learning, and information-theoretic security [3].

Data injection attacks (DIAs) have emerged as a major source of concern and exemplify the type of cybersecurity threats that specifically target power systems [4]. DIAs manipulate the state estimation process in the EMS by altering the measurements of the state variables without triggering the bad data detection mechanism put in place by the operator. In [4] it is shown that attacks that lie in the column space of the Jacobian measurement matrix are undetectable by testing the residual. To decrease the number of the sensors that need to be compromised by the attacker while remaining undetectable, the  $\ell_0$  norm of the attack vector is used as minimization objective yielding sparse attack in [5], [6], [7] and [8]. The case in which sparse attacks are constructed in a distributed setting with multiple attackers is discussed in [9] and [10].

The complex nature of the power system leads naturally to a stochastic modelling of the state variables describing the grid. For instance, the state variables of low voltage distribution systems are well described as following a multivariate Gaussian distribution [11]. DIAs within a Bayesian framework with minimum mean square error (MMSE) estimation are studied in [12] for the centralized case and in [13] for the distributed case. However, the fundamental limits governing the performance of attacks in the smart grid are not well understood yet.

Information-theoretic tools are well suited to analyze power system by leveraging the stochastic description of the state variables. A sensor placement strategy that accounts for the amount of information acquired by the sensing infrastructure is studied in [14]. Information-theoretic privacy guarantees for smart meter users are proposed in [15], [16], [17] for memoryless stochastic processes and in [18] for general random processes. In [19], stealth Gaussian DIA constructions are studied in terms of information measures that quantify the information loss and the probability of attack detection induced by the attack. Therein, the proposed cost function gives the same weight to the information loss and the probability of detection which results in the effective secrecy framework proposed by [20] in the context of *stealth* communications. Stealth DIA constructions are also studied in [5], [21] for the case in which the detection is based on the

residual and in a Bayesian hypothesis testing framework in [22]. The approaches in [5] and [21] consider the minimum cost of compromising the meters and the communication substation, respectively. On the other hand, [22] focuses on the delay between the time of attacker launching the attack and the time of operator detecting the attack.

In this paper, the stealth attacks in [19] are generalized by introducing a weight parameter to the objective describing the probability of detection, which allows the attacker to construct attacks with arbitrarily low probability of detection. Operating under the assumption that the state variables are described by a multivariate Gaussian distribution [12], [13], we characterize the optimal Gaussian generalized stealth attacks. Since the performance of the attacks depends on the weighting parameter governing the probability of detection, we provide a sufficient condition on the weighting parameter that achieves a desired probability of attack detection. To this end, we characterize the probability of attack detection via an upper bound which leverages a concentration inequality in [23].

The organization of the rest paper is shown as following: In Section II, a Bayesian framework with linearized dynamics for DIA is introduced. The generalized stealth attack construction and performance analysis are presented in Section III. Section IV provides the probability of detection of the generalized stealth attack, and the concentration inequality based upper bound for probability of detection. Section V verifies the results of Section III and Section IV on IEEE Test System. The paper ends with conclusions in Section VI.

## II. SYSTEM MODEL

### A. Bayesian Framework with Linearized Dynamics

The measurement model for state estimation with linearized dynamics is given by

$$Y^M = \mathbf{H}X^N + Z^M, \quad (1)$$

where  $Y^M \in \mathbb{R}^M$  is a vector of random variables describing the measurements;  $X^N \in \mathbb{R}^N$  is a vector of random variables describing the state variables;  $\mathbf{H} \in \mathbb{R}^{M \times N}$  is the linearized Jacobian measurement matrix which is determined by the power network topology and the admittances of the branches; and  $Z^M \in \mathbb{R}^M$  is the additive white Gaussian noise (AWGN) with distribution  $\mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_M)$  that is introduced by the sensors as a result of the thermal noise [24], [25]. In the

remaining of the paper, we assume that the vector of the state variables follows a multivariate Gaussian distribution given by

$$X^N \sim \mathcal{N}(\mathbf{0}, \Sigma_{XX}), \quad (2)$$

where  $\Sigma_{XX} \in \mathcal{S}_+^N$  is the covariance matrix of the distribution of the state variables and  $\mathcal{S}_+^N$  denotes the set of positive semidefinite matrices of size  $N \times N$ . As a result of the linearized dynamic in (1), the vector of measurements also follows a multivariate Gaussian distribution denoted by

$$Y^M \sim \mathcal{N}(\mathbf{0}, \Sigma_{YY}), \quad (3)$$

where  $\Sigma_{YY} = \mathbf{H}\Sigma_{XX}\mathbf{H}^T + \sigma^2\mathbf{I}_M$  is the covariance matrix of the distribution of the vector of measurements.

Data injection attacks corrupt the measurements available to the operator by adding an attack vector to the measurements. The resulting vector of compromised measurements is given by

$$Y_A^M = \mathbf{H}X^N + Z^M + A^M, \quad (4)$$

where  $A^M \in \mathbb{R}^M$  is the attack vector and  $Y_A^M \in \mathbb{R}^M$  is the vector containing the compromised measurements [4]. Given the stochastic nature of the state variables, it is reasonable for the attacker to pursue a stochastic attack construction strategy. In the following, an attack vector which is independent of the state variables is constructed under the assumption that the attack vector follows a multivariate Gaussian distribution denoted by

$$A^M \sim \mathcal{N}(\mathbf{0}, \Sigma_{AA}), \quad (5)$$

where  $\Sigma_{AA} \in \mathcal{S}_+^M$  is the covariance matrix of the attack distribution. The rationale for choosing a Gaussian distribution for the attack vector follows from the fact that for the measurement model in (4) the additive attack distribution that minimizes the mutual information between the vector of state variables and the compromised measurements is Gaussian [26]. Because of the Gaussianity of the attack distribution, the vector of compromised measurements is distributed as

$$Y_A^M \sim \mathcal{N}(\mathbf{0}, \Sigma_{Y_A Y_A}), \quad (6)$$

where  $\Sigma_{Y_A Y_A} = \mathbf{H}\Sigma_{XX}\mathbf{H}^T + \sigma^2\mathbf{I}_M + \Sigma_{AA}$  is the covariance matrix of the distribution of the compromised measurements.

It is worth noting that the independence of the attack vector with respect to the state variables implies that constructing the attack vector does not require access to the realizations of the state variables. In fact, knowledge of the second order moments of the state variables and the variance of the AWGN introduced by the measurement process suffices to construct the attack. This assumption significantly reduces the difficulty of the attack construction.

The operator of the power system makes use of the acquired measurements to detect the attack. The detection problem is cast as a hypothesis testing problem with hypotheses

$$\mathcal{H}_0 : Y^M \sim \mathcal{N}(\mathbf{0}, \Sigma_{YY}), \quad \text{versus} \quad (7)$$

$$\mathcal{H}_1 : Y^M \sim \mathcal{N}(\mathbf{0}, \Sigma_{Y_A Y_A}). \quad (8)$$

The null hypothesis  $\mathcal{H}_0$  describes the case in which the power system is not compromised, while the alternative hypothesis  $\mathcal{H}_1$  describes the case in which the power system is under attack.

Two types of error are considered in hypothesis testing problems, Type I error is the probability of a “true negative” event; and Type II error is the probability of a “false alarm” event. The Neyman-Pearson lemma [27] states that for a fixed probability of Type I error, the likelihood ratio test (LRT) achieves the minimum Type II error when compared with any other test with an equal or smaller Type I error. Consequently, the LRT is chosen to decide between  $\mathcal{H}_0$  and  $\mathcal{H}_1$  based on the available measurements. The LRT between  $\mathcal{H}_0$  and  $\mathcal{H}_1$  takes following form:

$$L(\mathbf{y}) \triangleq \frac{f_{Y_A^M}(\mathbf{y})}{f_{Y^M}(\mathbf{y})} \underset{\mathcal{H}_0}{\overset{\mathcal{H}_1}{\geq}} \tau, \quad (9)$$

where  $\mathbf{y} \in \mathbb{R}^M$  is a realization of the vector of random variables modelling the measurements,  $f_{Y_A^M}$  and  $f_{Y^M}$  denote the probability density functions (p.d.f.'s) of  $Y_A^M$  and  $Y^M$ , respectively, and  $\tau$  is the decision threshold set by the operator to meet the false alarm constraint.

### B. Information-Theoretic Setting

The mutual information between two random variables is a measure of the amount of information that each random variable contains about the other random variable. Consequently, the amount of information that the vector of measurements contains about the vector of state variables is determined by the mutual information between the vector of state variables and the vector of measurements. The Kullback-Leibler (KL) divergence between two probability distributions is a measure of the statistical similarity between the distributions. For the hypothesis testing problem in (9), a small value of the KL divergence between  $P_{Y_A^M}$  and  $P_{Y^M}$  implies that

on average the attack is unlikely to be detected by the LRT set by the attacker for a fixed value of  $\tau$ .

The purpose of the attacker is to disrupt the normal state estimation procedure by minimizing the information that the operator acquires about the state variables, while guaranteeing that the probability of attack detection is small enough, and therefore, remain concealed in the system.

An information-theoretic framework for the attack construction is adopted in this paper. To minimize the information that the operator acquires about the state variables from the measurements, the attacker minimizes the mutual information between the vector of state variables and the vector of compromised measurements. Specifically, the attacker aims to minimize  $I(X^N; Y_A^M)$ . On the other hand, the probability of attack detection is determined by the detection threshold  $\tau$  set by the operator and the distribution induced by the attack on the vector of compromised measurements. An analytical expression of the probability of attack detection can be described in closed-form as a function of the distributions describing the measurements under both hypotheses. However, the expression is involved in general and it is not straightforward to incorporate it into an analytical formulation of the attack construction. For that reason, we instead consider the asymptotic performance of the LRT to evaluate the detection performance of the operator. The Chernoff-Stein lemma [28] characterizes the asymptotic exponent of the probability of detection when the number of observations of measurement vectors grows to infinity. In our setting, the Chernoff-Stein lemma states that for any LRT and  $\epsilon \in (0, 1/2)$ , it holds that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^\epsilon = -D(P_{Y_A^M} || P_{Y^M}), \quad (10)$$

where  $D(\cdot || \cdot)$  is the KL divergence,  $\beta_n^\epsilon$  is the minimum Type II error such that the Type I error  $\alpha$  satisfies  $\alpha < \epsilon$ , and  $n$  is the number of  $M$ -dimensional measurement vectors that are available for the LRT. Therefore, for the attacker, minimizing the asymptotic detection probability is equivalent to minimizing  $D(P_{Y_A^M} || P_{Y^M})$ , where  $P_{Y_A^M}$  and  $P_{Y^M}$  denote the probability distributions of  $Y_A^M$  and  $Y^M$ , respectively.

### III. INFORMATION-THEORETIC ATTACK

#### A. Generalized Stealth Attacks

When these two information-theoretic objectives are considered by the attacker, [19] proposes an stealthy attack construction that combines the two objectives in one cost function, i.e.,

$$I(X^N; Y_A^M) + D(P_{Y_A^M} || P_{Y^M}) = D(P_{X^N Y_A^M} || P_{X^N} P_{Y^M}), \quad (11)$$

where  $P_{X^N Y_A^M}$  is the joint distribution of  $X^N$  and  $Y_A^M$ . The resulting optimization problem to construct the attack is given by

$$\min_{A^M} D(P_{X^N Y_A^M} || P_{X^N} P_{Y^M}). \quad (12)$$

Therein, it is shown that (12) is a convex optimization problem and the covariance matrix of the optimal Gaussian attack is  $\Sigma_{AA} = \mathbf{H}\Sigma_{XX}\mathbf{H}^T$ . However, numerical simulations on IEEE test system show that the attack construction proposed above yields large values of probability of detection in practical settings.

To address the issue of high probability of detection, in the following we propose an attack construction strategy that tunes the probability of detection with a parameter that weights the detection term in the cost function. The resulting optimization problem is given by

$$\min_{A^M} I(X^N; Y_A^M) + \lambda D(P_{Y_A^M} || P_{Y^M}), \quad (13)$$

where  $\lambda \geq 1$  governs the weight given to each objective in the cost function. It is interesting to note that for the case in which  $\lambda = 1$  the proposed cost function boils down to the effective secrecy proposed in [20] and the attack construction in (13) coincides with that in [19]. For  $\lambda > 1$ , the attacker adopts a conservative approach and prioritizes remaining undetected over minimizing the amount of information acquired by the operator. By increasing the value of  $\lambda$  the attacker decreases the probability of detection at the expense of increasing the amount of information acquired by the operator via the measurements.

### B. Optimal Attack Construction

The attack construction in (13) is formulated in a general setting. The following propositions particularize the KL divergence and MI to our multivariate Gaussian setting.

**Proposition 1.** [28] *The KL divergence between  $M$ -dimensional multivariate Gaussian distributions  $\mathcal{N}(\mathbf{0}, \Sigma_{Y_A Y_A})$  and  $\mathcal{N}(\mathbf{0}, \Sigma_{YY})$  is given by*

$$D(P_{Y_A^M} || P_{Y^M}) = \frac{1}{2} \left( \log \frac{|\Sigma_{YY}|}{|\Sigma_{Y_A Y_A}|} - M + \text{tr}(\Sigma_{YY}^{-1} \Sigma_{Y_A Y_A}) \right). \quad (14)$$

**Proposition 2.** [28] *The mutual information between the vectors of random variables  $X^N \sim \mathcal{N}(\mathbf{0}, \Sigma_{XX})$  and  $Y_A^M \sim \mathcal{N}(\mathbf{0}, \Sigma_{Y_A Y_A})$  is given by*

$$I(X^N; Y_A^M) = \frac{1}{2} \log \frac{|\Sigma_{XX}| |\Sigma_{Y_A Y_A}|}{|\Sigma|}, \quad (15)$$



where  $\Sigma$  is the covariance matrix of the joint distribution of  $(X^N, Y_A^M)$ .

Substituting (14) and (15) in (13) we can now pose the Gaussian attack construction as the following optimization problem:

$$\begin{aligned} \min_{\Sigma_{AA} \in \mathcal{S}_+^M} & -(\lambda - 1) \log |\Sigma_{YY} + \Sigma_{AA}| - \log |\Sigma_{AA} + \sigma^2 \mathbf{I}_M| \\ & + \lambda \text{tr}(\Sigma_{YY}^{-1} \Sigma_{AA}). \end{aligned} \quad (16)$$

We now proceed to solve the optimization problem above. First, note that the optimization domain  $\mathcal{S}_+^M$  is a convex set. The following proposition characterizes the convexity of the cost function.

**Proposition 3.** *Let  $\lambda \geq 1$ . Then the cost function in the optimization problem in (16) is convex.*

*Proof.* Note that the term  $-\log |\Sigma_{AA} + \sigma^2 \mathbf{I}_M|$  is a convex function on  $\Sigma_{AA} \in \mathcal{S}_+^M$  [29]. Additionally,  $-(\lambda - 1) \log |\Sigma_{YY} + \Sigma_{AA}|$  is a convex function on  $\Sigma_{AA} \in \mathcal{S}_+^M$  when  $\lambda \geq 1$ . Since the trace operator is a linear operator and the sum of convex functions is convex, it follows that the cost function in (16) is convex on  $\Sigma_{AA} \in \mathcal{S}_+^M$ .  $\square$

**Theorem 1.** *Let  $\lambda \geq 1$ . Then the solution to the optimization problem in (16) is*

$$\Sigma_{AA}^* = \frac{1}{\lambda} \mathbf{H} \Sigma_{XX} \mathbf{H}^T. \quad (17)$$

*Proof.* Denote the cost function in (16) by  $f(\Sigma_{AA})$ . Taking the derivative of the cost function with respect to  $\Sigma_{AA}$  yields

$$\begin{aligned} \frac{\partial f(\Sigma_{AA})}{\partial \Sigma_{AA}} &= -2(\lambda - 1)(\Sigma_{YY} + \Sigma_{AA})^{-1} - 2(\Sigma_{AA} + \sigma^2 \mathbf{I}_M)^{-1} \\ &\quad + 2\lambda \Sigma_{YY}^{-1} + (\lambda - 1) \text{diag}((\Sigma_{YY} + \Sigma_{AA})^{-1}) \\ &\quad + \text{diag}((\Sigma_{AA} + \sigma^2 \mathbf{I}_M)^{-1}) - \lambda \text{diag}(\Sigma_{YY}^{-1}). \end{aligned} \quad (18)$$

Note that the only critical point is  $\Sigma_{AA}^* = \frac{1}{\lambda} \mathbf{H} \Sigma_{XX} \mathbf{H}^T$ . Theorem 1 follows immediately from combining this result with Proposition 3.  $\square$

**Corollary 1.** *The mutual information between the vector of state variables and the vector of compromised measurements induced by the optimal attack construction is given by*

$$\begin{aligned} I(X^N; Y_A^M) &= \frac{1}{2} \log \left| \mathbf{H} \Sigma_{XX} \mathbf{H}^T \left( \sigma^2 \mathbf{I}_M + \frac{1}{\lambda} \mathbf{H} \Sigma_{XX} \mathbf{H}^T \right)^{-1} + \mathbf{I}_M \right|. \end{aligned} \quad (19)$$

Theorem 1 shows that the generalized stealth attacks share the same structure of the stealth attacks in [19] up to a scaling factor determined by  $\lambda$ . The solution in Theorem 1 holds for the case in which  $\lambda \geq 1$ , and therefore, lacks full generality. However, the case in which  $\lambda < 1$  yields unreasonably high probability of detection [19] which indicates that the proposed attack construction is indeed of practical interest in a wide range of state estimation settings.

The resulting attack construction is remarkably simple to implement provided that the information about the system is available to the attacker. Indeed, the attacker only requires access to the linearized Jacobian measurement matrix  $\mathbf{H}$  and the second order statistics of the state variables, but the variance of the noise introduced by the sensors is not necessary. To obtain the Jacobian, a malicious attacker needs to know the topology of the grid, the admittances of the branches, and the operation point of the system. The second order statistics of the state variables on the other hand, can be estimated using historical data. In [19] it is shown that the attack construction with a sample covariance matrix of the state variables obtained with historical data is asymptotically optimal when the size of the training data grows to infinity.

It is interesting to note that the mutual information increases monotonically with  $\lambda$  and that it asymptotically converges to  $I(X^N; Y^M)$ , i.e. the case in which there is no attack. While the evaluation of the mutual information as shown in Corollary 1 is straightforward, the computation of the associated probability of detection yields involved expressions that do not provide much insight. For that reason, the probability of detection of optimal attacks is treated in the following section.

#### IV. PROBABILITY OF DETECTION OF GENERALIZED STEALTH ATTACKS

The asymptotic probability of detection of the generalized stealth attacks characterized in Section III-B is governed by the KL divergence as described in (10). However in the non-asymptotic case, determining the probability of detection is difficult, and therefore, choosing a value of  $\lambda$  that provides the desired probability of detection is a challenging task. In this section we first provide a closed-form expression of the probability of detection by direct evaluation and show that the expression does not provide any practical insight over the choice of  $\lambda$  that achieves the desired detection performance. That being the case, we then provide an upper bound on the probability of detection, which, in turn, provides a lower bound on the value of  $\lambda$  that achieves the desired probability of detection.

### A. Direct Evaluation of the Probability of Detection

Detection based on the LRT with threshold  $\tau$  yields a probability of detection given by

$$P_D \triangleq \mathbb{E} \left[ \mathbb{1}_{\{L(Y_A^M) \geq \tau\}} \right], \quad (20)$$

where  $\mathbb{1}_{\{\cdot\}}$  is the indicator function. The following proposition particularizes the above expression to the optimal attack construction described in Section III-B.

**Lemma 1.** *The probability of detection of the LRT in (9) for the attack construction in (17) is given by*

$$P_D(\lambda) = \mathbb{P} \left[ (U^P)^T \Delta U^P \geq \lambda (2 \log \tau + \log |\mathbf{I}_P + \lambda^{-1} \Delta|) \right], \quad (21)$$

where  $P = \text{rank}(\mathbf{H}\Sigma_{XX}\mathbf{H}^T)$ ,  $U^P \in \mathbb{R}^P$  is a vector of random variables with distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I}_P)$ , and  $\Delta \in \mathbb{R}^{P \times P}$  is a diagonal matrix with entries given by  $(\Delta)_{i,i} = \lambda_i(\mathbf{H}\Sigma_{XX}\mathbf{H}^T)\lambda_i(\Sigma_{YY}^{-1})$ , where  $\lambda_i(\mathbf{A})$  with  $i = 1, \dots, P$  denotes the  $i$ -th eigenvalue of matrix  $\mathbf{A}$  in descending order.

*Proof.* The probability of detection of the stealth attack is,

$$P_D(\lambda) = \int_{\mathcal{S}} dP_{Y_A^M} \quad (22)$$

$$= \frac{1}{(2\pi)^{\frac{M}{2}} |\Sigma_{Y_A Y_A}|^{\frac{1}{2}}} \int_{\mathcal{S}} \exp \left\{ -\frac{1}{2} \mathbf{y}^T \Sigma_{Y_A Y_A}^{-1} \mathbf{y} \right\} d\mathbf{y}, \quad (23)$$

where

$$\mathcal{S} = \{\mathbf{y} \in \mathbb{R}^M : L(\mathbf{y}) \geq \tau\}. \quad (24)$$

Algebraic manipulation yields the following equivalent description of the integration domain:

$$\mathcal{S} = \{\mathbf{y} \in \mathbb{R}^M : \mathbf{y}^T \Delta_0 \mathbf{y} \geq 2 \log \tau + \log |\mathbf{I}_M + \Sigma_{AA} \Sigma_{YY}^{-1}| \}, \quad (25)$$

with  $\Delta_0 \triangleq \Sigma_{YY}^{-1} - \Sigma_{Y_A Y_A}^{-1}$ . Let  $\Sigma_{YY} = \mathbf{U}_{YY} \Lambda_{YY} \mathbf{U}_{YY}^T$  where  $\Lambda_{YY} \in \mathbb{R}^{M \times M}$  is a diagonal matrix containing the eigenvalues of  $\Sigma_{YY}$  in descending order and  $\mathbf{U}_{YY} \in \mathbb{R}^{M \times M}$  is a unitary matrix whose columns are the eigenvectors of  $\Sigma_{YY}$  ordered matching the order of the eigenvalues. Applying the change of variable  $\mathbf{y}_1 \triangleq \mathbf{U}_{YY} \mathbf{y}$  in (23) results in

$$P_D(\lambda) = \frac{1}{(2\pi)^{\frac{M}{2}} |\Sigma_{Y_A Y_A}|^{\frac{1}{2}}} \int_{\mathcal{S}_1} \exp \left\{ -\frac{1}{2} \mathbf{y}_1^T \Lambda_{Y_A Y_A}^{-1} \mathbf{y}_1 \right\} d\mathbf{y}_1, \quad (26)$$

where  $\Lambda_{Y_A Y_A} \in \mathbb{R}^{M \times M}$  denotes the diagonal matrix containing the eigenvalues of  $\Sigma_{Y_A Y_A}$  in descending order. Noticing that  $\Sigma_{YY}$ ,  $\Sigma_{AA}$  and  $\Sigma_{Y_A Y_A}$  are also diagonalized by  $\mathbf{U}_{YY}$ , the integration domain  $\mathcal{S}_1$  is given by

$$\mathcal{S}_1 = \{\mathbf{y}_1 \in \mathbb{R}^M : \mathbf{y}_1^T \Delta_1 \mathbf{y}_1 \geq 2 \log \tau + \log |\mathbf{I}_M + \Lambda_{AA} \Lambda_{YY}^{-1}| \}, \quad (27)$$

where  $\Delta_1 \triangleq \Lambda_{YY}^{-1} - \Lambda_{Y_A Y_A}^{-1}$  with  $\Lambda_{AA}$  denoting the diagonal matrix containing the eigenvalues of  $\Sigma_{AA}$  in descending order. Further applying the change of variable  $\mathbf{y}_2 \triangleq \Lambda_{Y_A Y_A}^{-\frac{1}{2}} \mathbf{y}_1$  in (26) results in

$$P_D(\lambda) = \frac{1}{\sqrt{(2\pi)^M}} \int_{\mathcal{S}_2} \exp\left\{-\frac{1}{2} \mathbf{y}_2^T \mathbf{y}_2\right\} d\mathbf{y}_2, \quad (28)$$

with the transformed integration domain given by

$$\mathcal{S}_2 = \{\mathbf{y}_2 \in \mathbb{R}^M: \mathbf{y}_2^T \Delta_2 \mathbf{y}_2 \geq 2 \log \tau + \log |\mathbf{I}_M + \Delta_2|\}, \quad (29)$$

with

$$\Delta_2 \triangleq \Lambda_{AA} \Lambda_{YY}^{-1}. \quad (30)$$

Setting  $\Delta \triangleq \lambda \Delta_2$  and noticing that  $\text{rank}(\Delta) = \text{rank}(\mathbf{H} \Sigma_{XX} \mathbf{H}^T)$  concludes the proof.  $\square$

Notice that the left-hand term  $(U^M)^T \Delta U^M$  in (21) is a weighted sum of independent  $\chi^2$  distributed random variables with one degree of freedom where the weights are determined by the diagonal entries of  $\Delta$  which depend on the second order statistics of the state variables, the Jacobian measurement matrix, and the variance of the noise; i.e. the attacker has no control over this term. The right-hand side contains in addition  $\lambda$  and  $\tau$ , and therefore, the probability of attack detection is described as a function of the parameter  $\lambda$ . However, characterizing the distribution of the resulting random variable is not practical since there is no closed-form expression for the distribution of a positively weighted sum of independent  $\chi^2$  random variables with one degree of freedom [30]. Usually, some moment matching approximation approaches such as the Lindsay–Pilla–Basak (LPB) method [31] are utilized to solve this problem but the resulting expressions are complex and the relation of the probability of detection with  $\lambda$  is difficult to describe analytically following this course of action. In the following an upper bound on the probability of attack detection is derived. The upper bound is then used to provide a simple lower bound on the value  $\lambda$  that achieves the desired probability of detection.

### B. Upper Bound on the Probability of Detection

The following theorem provides a sufficient condition for  $\lambda$  to achieve a desired probability of attack detection.

**Theorem 2.** *Let  $\tau > 1$  be the decision threshold of the LRT. For any  $t > 0$  and  $\lambda \geq \max(\lambda^*(t), 1)$  then the probability of attack detection satisfies*

$$P_D(\lambda) \leq e^{-t}, \quad (31)$$

where  $\lambda^*(t)$  is the only positive solution of  $\lambda$  satisfying

$$2\lambda \log \tau - \frac{1}{2\lambda} \text{tr}(\Delta^2) - 2\sqrt{\text{tr}(\Delta^2)t} - 2\|\Delta\|_\infty t = 0. \quad (32)$$

*Proof.* We start with the result of Lemma 1 which gives

$$P_D(\lambda) = \mathbb{P} \left[ (U^P)^T \Delta U^P \geq \lambda (2 \log \tau + \log |\mathbf{I}_P + \lambda^{-1} \Delta|) \right]. \quad (33)$$

We now proceed to expand the term  $\log |\mathbf{I}_P + \lambda^{-1} \Delta|$  using a Taylor series expansion resulting in

$$\begin{aligned} \log |\mathbf{I}_P + \lambda^{-1} \Delta| &= \sum_{i=1}^P \log (1 + \lambda^{-1} (\Delta)_{i,i}) \\ &= \sum_{i=1}^P \left( \sum_{j=1}^{\infty} \left( \frac{(\lambda^{-1} (\Delta)_{i,i})^{2j-1}}{2j-1} - \frac{(\lambda^{-1} (\Delta)_{i,i})^{2j}}{2j} \right) \right). \end{aligned} \quad (34)$$

$$= \sum_{i=1}^P \left( \sum_{j=1}^{\infty} \left( \frac{(\lambda^{-1} (\Delta)_{i,i})^{2j-1}}{2j-1} - \frac{(\lambda^{-1} (\Delta)_{i,i})^{2j}}{2j} \right) \right). \quad (35)$$

Since  $(\Delta)_{i,i} \leq 1$ , for  $i = 1, \dots, P$ , and  $\lambda \geq 1$ , then

$$\frac{(\lambda^{-1} (\Delta)_{i,i})^{2j-1}}{2j-1} - \frac{(\lambda^{-1} (\Delta)_{i,i})^{2j}}{2j} \geq 0, \text{ for } j \in \mathbb{Z}^+. \quad (36)$$

Thus, (35) is lower bounded by the second order Taylor expansion, i.e.,

$$\log |\mathbf{I}_P + \Delta| \geq \sum_{i=1}^P \left( \lambda^{-1} (\Delta)_{i,i} - \frac{(\lambda^{-1} (\Delta)_{i,i})^2}{2} \right) \quad (37)$$

$$= \frac{1}{\lambda} \text{tr}(\Delta) - \frac{1}{2\lambda^2} \text{tr}(\Delta^2). \quad (38)$$

Substituting (38) in (33) yields

$$P_D(\lambda) \leq \mathbb{P} \left[ (U^P)^T \Delta U^P \geq \text{tr}(\Delta) + 2\lambda \log \tau - \frac{1}{2\lambda} \text{tr}(\Delta^2) \right]. \quad (39)$$

Note that  $\mathbb{E} [(U^P)^T \Delta U^P] = \text{tr}(\Delta)$ , and therefore, evaluating the probability in (39) is equivalent to evaluating the probability of  $(U^P)^T \Delta U^P$  deviating  $2\lambda \log \tau - \frac{1}{2\lambda} \text{tr}(\Delta^2)$  from the mean. In view of this, the right-hand side in (39) is upper bounded by [23], [32]

$$P_D(\lambda) \leq \mathbb{P} \left[ (U^P)^T \Delta U^P \geq \text{tr}(\Delta) + 2\sqrt{\text{tr}(\Delta^2)t} + 2\|\Delta\|_\infty t \right] \quad (40)$$

$$\leq e^{-t}, \quad (41)$$

for  $t > 0$  satisfying

$$2\lambda \log \tau - \frac{1}{2\lambda} \text{tr}(\Delta^2) \geq 2\sqrt{\text{tr}(\Delta^2)t} + 2\|\Delta\|_\infty t. \quad (42)$$

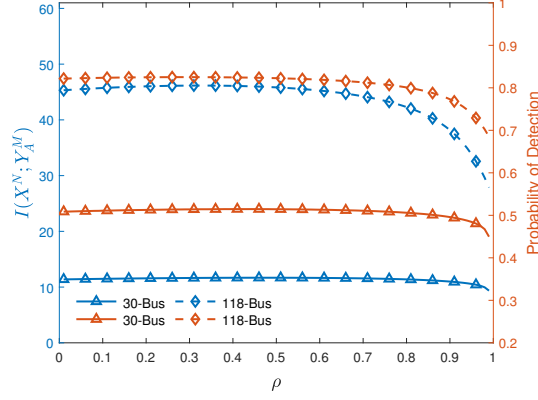


Fig. 1. Performance of the generalized stealth attack in terms of mutual information and probability of detection for different values of  $\rho$  when  $\lambda = 2$ ,  $\tau = 2$ , and SNR = 10 dB.

The expression in (42) is satisfied with equality for two values of  $\lambda$ , one is strictly negative and the other one is strictly positive denoted by  $\lambda^*(t)$ , when  $\tau > 1$ . The result follows by noticing that the left-hand term of (42) increases monotonically for  $\lambda > 0$  and choosing  $\lambda \geq \max(\lambda^*(t), 1)$ . This concludes the proof.  $\square$

It is interesting to note that for large values of  $\lambda$  the probability of detection decreases exponentially fast with  $\lambda$ . We will later show in the numerical results that the regime in which the exponentially fast decrease kicks in does not align with the saturation of the mutual information loss induced by the attack.

## V. NUMERICAL SIMULATION

In this section, we present simulations to evaluate the performance of the proposed attack strategy in practical state estimation settings. In particular, the IEEE 30-Bus and 118-Bus test systems are utilized in the simulation. In state estimation with linearized dynamics, the Jacobian measurement matrix is determined by the operation point. We assume a DC state estimation scenario [24], [25], and thus, we set the bus voltage angles to zero. Note that in this setting it is sufficient to specify the network topology, the branch reactances, real power flow, and the power injection values to fully characterize the system. Specifically, we use the IEEE test system framework provided by MATPOWER [33].

As stated in Section IV-A, there is no closed-form expression for the distribution of a positively weighted sum of independent  $\chi^2$  random variables, which is required to calculate the probability

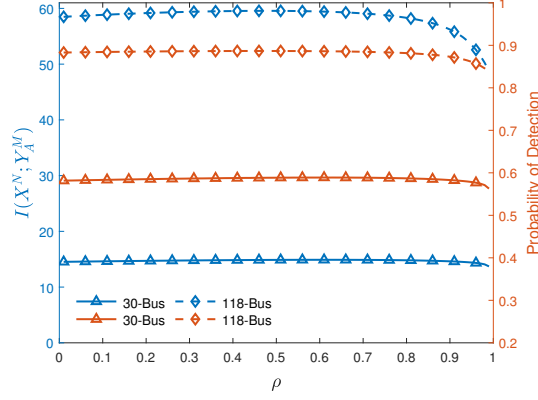


Fig. 2. Performance of the generalized stealth attack in terms of mutual information and probability of detection for different values of  $\rho$  when  $\lambda = 2$ ,  $\tau = 2$ , and SNR = 20 dB.

of detection of the generalized stealth attacks as shown in Lemma 1. For that reason, we use the LPB method and the MOMENTCHI2 package [34] to numerically evaluate the probability of attack detection.

The simulation settings are the same as in [19]. The covariance matrix of the state variables is assumed to be a Toeplitz matrix with exponential decay parameter  $\rho$ , where the exponential decay parameter  $\rho$  determines the correlation strength between different entries of the state variable vector. The performance of the generalized stealth attack is a function of weight given to the detection term in the attack construction cost function, i.e.  $\lambda$ , the correlation strength between state variables, i.e.  $\rho$ , and the Signal-to-Noise Ratio (SNR) of the power system which is defined as

$$\text{SNR} \triangleq 10 \log_{10} \left( \frac{\text{tr}(\mathbf{H}\mathbf{\Sigma}_{XX}\mathbf{H}^T)}{M\sigma^2} \right). \quad (43)$$

Fig. 1 and Fig. 2 depict the performance of the optimal attack construction given in (17) for different values of  $\rho$  with SNR = 10 dB and SNR = 20 dB, respectively, when  $\lambda = 2$  and  $\tau = 2$ . Interestingly, the performance of the attack construction does not change monotonically with the correlation strength, which suggests that the correlation among the state variables does not necessarily provide an advantage to the attacker. Admittedly, for a small and moderate values of  $\rho$ , the performance of the attack does not change significantly with  $\rho$  for both objectives. This effect is more noticeable in the high SNR scenario. However, for large values of  $\rho$  the performance of the attack improves significantly in terms of both mutual information and probability of detection. Moreover, the advantage provided by large values of  $\rho$  is more significant for the

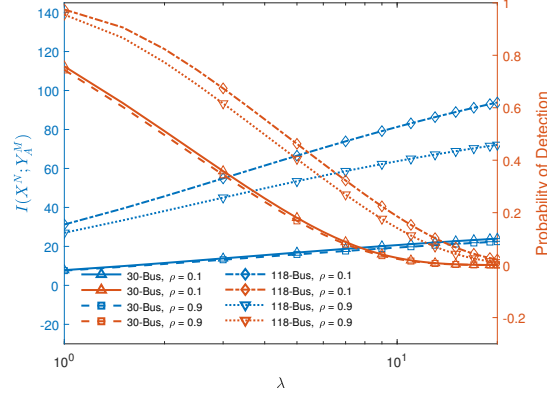


Fig. 3. Performance of the generalized stealth attack in terms of mutual information and probability of detection for different values of  $\lambda$  and system size when  $\rho = 0.1$ ,  $\rho = 0.9$ , SNR = 10 dB and  $\tau = 2$ .

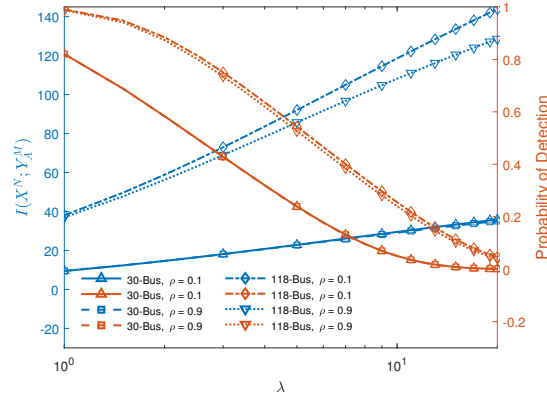


Fig. 4. Performance of the generalized stealth attack in terms of mutual information and probability of detection for different values of  $\lambda$  and system size when  $\rho = 0.1$ ,  $\rho = 0.9$ , SNR = 20 dB and  $\tau = 2$ .

118-Bus system than for the 30-Bus system, which indicates that correlation between the state variables is easier to exploit for the attacker in large systems.

Fig. 3 and Fig. 4 depict the performance of the optimal attack construction for different values of  $\lambda$  and  $\rho$  with SNR = 10 dB and SNR = 20 dB, respectively, when  $\tau = 2$ . As expected, larger values of the parameter  $\lambda$  yield smaller values of the probability of attack detection while increasing the mutual information between the state variables vector and the compromised measurement vector. We observe that the probability of detection decreases approximately linearly for moderate values of  $\lambda$ . On the other hand, Theorem 2 states that for large values of  $\lambda$  the



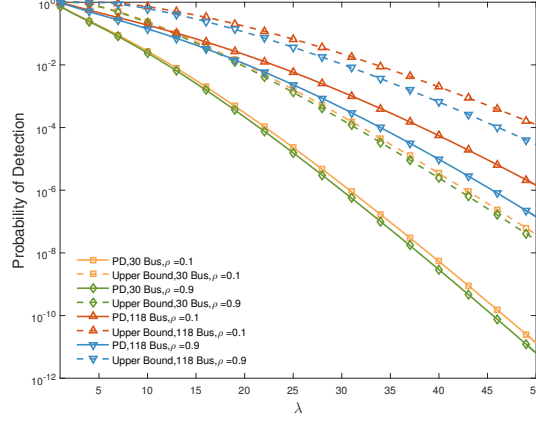


Fig. 5. Upper bound on probability of detection given in Theorem 2 for different values of  $\lambda$  when  $\rho = 0.1$  or  $0.9$ ,  $\text{SNR} = 10$  dB, and  $\tau = 2$ .

probability of detection decreases exponentially fast to zero. However, for the range of values of  $\lambda$  in which the decrease of probability of detection is approximately linear, there is no significant reduction on the rate of growth of mutual information. In view of this, the attacker needs to choose the value of  $\lambda$  carefully as the convergence of the mutual information to the asymptote  $I(X^N; Y^M)$  is slower than that of the probability of detection to zero.

The comparison between the 30-Bus and 118-Bus systems shows that for the smaller size system the probability of detection decreases faster to zero while the rate of growth of mutual information is smaller than that on the larger system. This suggests that the choice of  $\lambda$  is particularly critical in large size systems as smaller size systems exhibit a more robust attack performance for different values of  $\lambda$ . The effect of the correlation between the state variables is significantly more noticeable for the 118-bus system. While there is a performance gain for the 30-bus system in terms of both mutual information and probability of detection due to the high correlation between the state variables, the improvement is more noteworthy for the 118-bus case. Remarkably, the difference in terms of mutual information between the case in which  $\rho = 0.1$  and  $\rho = 0.9$  increases as  $\lambda$  increases which indicates that the cost in terms of mutual information of reducing the probability of detection is large in the small values of correlation.

The performance of the upper bound given by Theorem 2 on the probability of detection for different values of  $\lambda$  and  $\rho$  when  $\tau = 2$  and  $\text{SNR} = 10$  dB is shown in Fig. 5. Similarly, Fig. 6 depicts the upper bound with the same parameters but with  $\text{SNR} = 20$  dB. As shown by Theorem 2 the bound decreases exponentially fast for large values of  $\lambda$ . Still, there is a significant gap

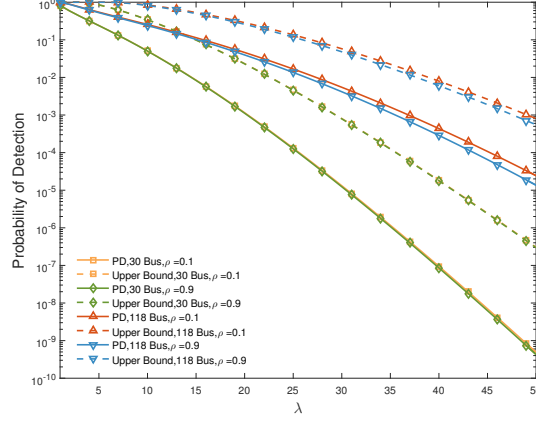


Fig. 6. Upper bound on probability of detection given in Theorem 2 for different values of  $\lambda$  when  $\rho = 0.1$  or  $0.9$ ,  $\text{SNR} = 20$  dB, and  $\tau = 2$ .

to the probability of attack detection evaluated numerically. This is partially due to the fact that our bound is based on the concentration inequality in [23] which introduces a gap of more than an order of magnitude. Interestingly, the gap decreases when the value of  $\rho$  increases although the change is not significant. More importantly, the bound is tighter for lower values of SNR for both 30-bus and 118-bus systems.

## VI. CONCLUSIONS

We have proposed a novel data injection attacks based on information-theoretic performance measures. Specifically, we have posed the attack construction problem as an optimization problem in which the cost function combines the mutual information and the probability of attack detection. The proposed cost function allows to obtain an arbitrarily small probability of attack detection via a parameter that weights the effect of the mutual information and the probability of detection. The resulting random attack construction has been analyzed in terms of the information loss and the probability of attack detection that it induces on the system. We have characterized the probability of attack detection by obtaining an easy to compute upper bound. The upper bound has been used to provide a practical attack construction guideline by determining the cost function that achieves a given probability of attack detection.

## REFERENCES

- [1] U.S.-Canada Power System Outage Task Force, *Final Report on the August 14, 2003 Blackout in the United States and Canada: Causes and Recommendations*, 2004.

- [2] D. Alderson and R. Di Pietro, "Operational technology: Are you vulnerable?," *Governance Directions*, vol. 68, no. 6, pp. 339–343, Jul. 2016.
- [3] Royal Society, *Progress and Research in Cybersecurity: Supporting a Resilient and Trustworthy System for the UK*, 2016.
- [4] Y. Liu, P. Ning, and M. K. Reiter, "False data injection attacks against state estimation in electric power grids," in *Proc. ACM Conf. on Computer and Communications Security*, Chicago, IL, USA, Nov. 2009, pp. 21–32.
- [5] G. Dan and H. Sandberg, "Stealth attacks and protection schemes for state estimators in power systems," in *Proc. IEEE Int. Conf. on Smart Grid Comm.*, Gaithersburg, MD, USA, Oct. 2010, pp. 214–219.
- [6] H. Sandberg, A. Teixeira, and K. H. Johansson, "On security indices for state estimators in power networks," in *Proc. 1st Workshop on Secure Control Systems*, Stockholm, Sweden, Apr. 2010.
- [7] K. C. Sou, H. Sandberg, and K. H. Johansson, "On the exact solution to a smart grid cyber-security analysis problem," *IEEE Trans. Smart Grid*, vol. 4, no. 2, pp. 856–865, Jun. 2013.
- [8] T. T. Kim and H. V. Poor, "Strategic protection against data injection attacks on power grids," *IEEE Trans. Smart Grid*, vol. 2, no. 2, pp. 326–333, Jun. 2011.
- [9] A. Tajer, S. Kar, H. V. Poor, and S. Cui, "Distributed joint cyber attack detection and state recovery in smart grids," in *Proc. IEEE Int. Conf. on Smart Grid Commun.*, Brussels, Belgium, Oct. 2011, pp. 202–207.
- [10] M. Ozay, I. Esnaola, F. T. Y. Vural, S. R. Kulkarni, and H. V. Poor, "Sparse attack construction and state estimation in the smart grid: Centralized and distributed models," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 7, pp. 1306–1318, Jul. 2013.
- [11] C. Genes, I. Esnaola, S. M. Perlaza, L. F. Ochoa, and D. Coca, "Recovering missing data via matrix completion in electricity distribution systems," in *Proc. IEEE Workshop on Signal Processing Advances in Wireless Communications*, Edinburgh, UK, Jul. 2016, pp. 1–6.
- [12] O. Kosut, L. Jia, R. J. Thomas, and L. Tong, "Malicious data attacks on the smart grid," *IEEE Trans. Smart Grid*, vol. 2, no. 4, pp. 645–658, Dec. 2011.
- [13] I. Esnaola, S. M. Perlaza, H. V. Poor, and O. Kosut, "Maximum distortion attacks in electricity grids," *IEEE Trans. Smart Grid*, vol. 7, no. 4, pp. 2007–2015, Jul. 2016.
- [14] Q. Li, T. Cui, Y. Weng, R. Negi, F. Franchetti, and M. D. Ilić, "An information-theoretic approach to PMU placement in electric power systems," *IEEE Trans. Smart Grid*, vol. 4, no. 1, pp. 446–456, Mar. 2013.
- [15] D. Varodayan and A. Khisti, "Smart meter privacy using a rechargeable battery: Minimizing the rate of information leakage," in *IEEE Int. Conf. on Acoust., Speech and Signal Process.*, Prague, Czech Republic, May 2011, pp. 1932–1935.
- [16] L. Sankar, S.R. Rajagopalan, S. Mohajer, and H.V. Poor, "Smart meter privacy: A theoretical framework," *IEEE Trans. Smart Grid*, vol. 4, no. 2, pp. 837–846, Jun. 2013.
- [17] O. Tan, D. Gündüz, and H. V. Poor, "Increasing smart meter privacy through energy harvesting and storage devices," *IEEE J. Sel. Areas Commun.*, vol. 31, no. 7, pp. 1331–1341, Jul. 2013.
- [18] M. Arrieta and I. Esnaola, "Smart meter privacy via the trapdoor channel," in *Proc. IEEE Int. Conf. on Smart Grid Commun.*, Dresden, Germany, Oct. 2017, pp. 227–282.
- [19] K. Sun, I. Esnaola, S.M. Perlaza, and H.V. Poor, "Information-theoretic attacks in the smart grid," in *Proc. IEEE Int. Conf. on Smart Grid Comm.*, Dresden, Germany, Oct. 2017, pp. 455–460.
- [20] J. Hou and G. Kramer, "Effective secrecy: Reliability, confusion and stealth," in *Proc. IEEE Int. Symp. on Information Theory*, Honolulu, HI, USA, Jun. 2014, pp. 601–605.
- [21] O. Vuković, K. C. Sou, G. Dán, and H. Sandberg, "Network-layer protection schemes against stealth attacks on state estimators in power systems," in *IEEE Int. Conf. on Smart Grid Comm.*, Brussels, Belgium, Oct. 2011, pp. 184–189.
- [22] Y. Huang, M. Esmalifalak, H. Nguyen, R. Zheng, Z. Han, H. Li, and L. Song, "Bad data injection in smart grid: Attack and defense mechanisms," *IEEE Commun. Mag.*, vol. 51, no. 1, pp. 27–33, Jan. 2013.

- [23] B. Laurent and P. Massart, “Adaptive estimation of a quadratic functional by model selection,” *Ann. Statist.*, vol. 28, no. 5, pp. 1302–1338, 2000.
- [24] A. Abur and A. G. Expósito, *Power System State Estimation: Theory and Implementation*, CRC Press, Mar. 2004.
- [25] J. J. Grainger and W. D. Stevenson, *Power System Analysis*, McGraw-Hill, 1994.
- [26] I. Shomorony and A. S. Avestimehr, “Worst-case additive noise in wireless networks,” *IEEE Trans. Inf. Theory*, vol. 59, no. 6, pp. 3833–3847, Jun. 2013.
- [27] J. Neyman and E. S. Pearson, “On the problem of the most efficient tests of statistical hypotheses,” in *Breakthroughs in Statistics*, Springer Series in Statistics, pp. 73–108. Springer New York, 1992.
- [28] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, John Wiley & Sons, Nov. 2012.
- [29] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, Mar. 2004.
- [30] D. A. Bodenham and N. M. Adams, “A comparison of efficient approximations for a weighted sum of chi-squared random variables,” *Stat Comput*, vol. 26, no. 4, pp. 917–928, Jul. 2016.
- [31] Bruce G. Lindsay, Ramani S. Pilla, and Prasanta Basak, “Moment-based approximations of distributions using mixtures: Theory and applications,” *Ann. Inst. Stat. Math.*, vol. 52, no. 2, pp. 215–230, Jun. 2000.
- [32] D. Hsu, S.M. Kakade, and T. Zhang, “A tail inequality for quadratic forms of subgaussian random vectors,” *Electron. Commun. in Probab.*, vol. 17, no. 52, pp. 1–6, 2012.
- [33] R. D. Zimmerman, C. E. Murillo-Sánchez, and R. J. Thomas, “MATPOWER: Steady-state operations, planning, and analysis tools for power systems research and education,” *IEEE Trans. Power Syst.*, vol. 26, no. 1, pp. 12–19, Feb. 2011.
- [34] D. Bodenham, *Momentchi2: Moment-Matching Methods for Weighted Sums of Chi-Squared Random Variables*. (2016) [Online]. Available: <https://cran.r-project.org/web/packages/momentchi2/index.html>.